# Continually Improving Grounded Natural Language Understanding through Human-Robot Dialog

**Jesse Thomason**
University of Texas at Austin
Ph.D. Defense
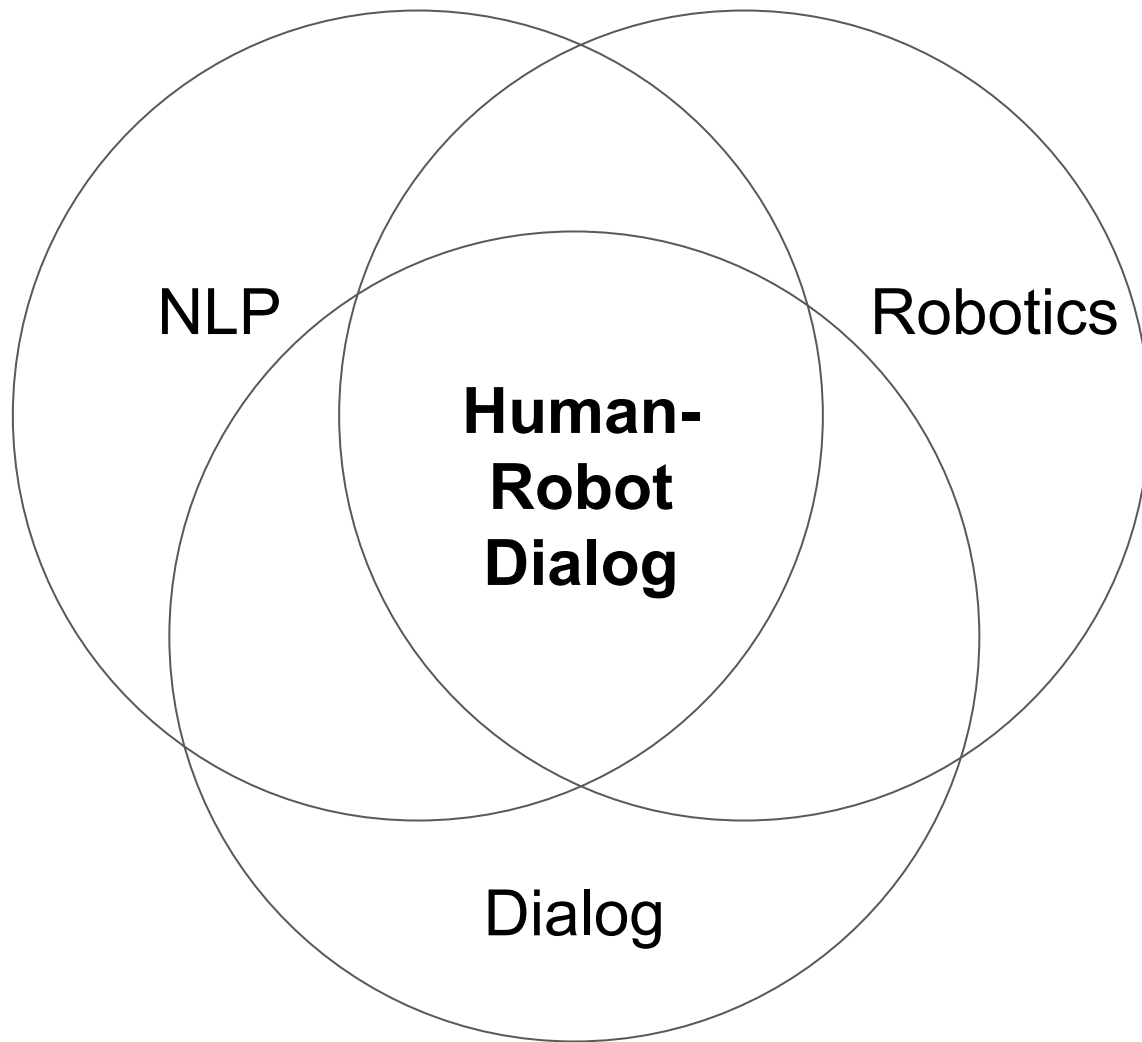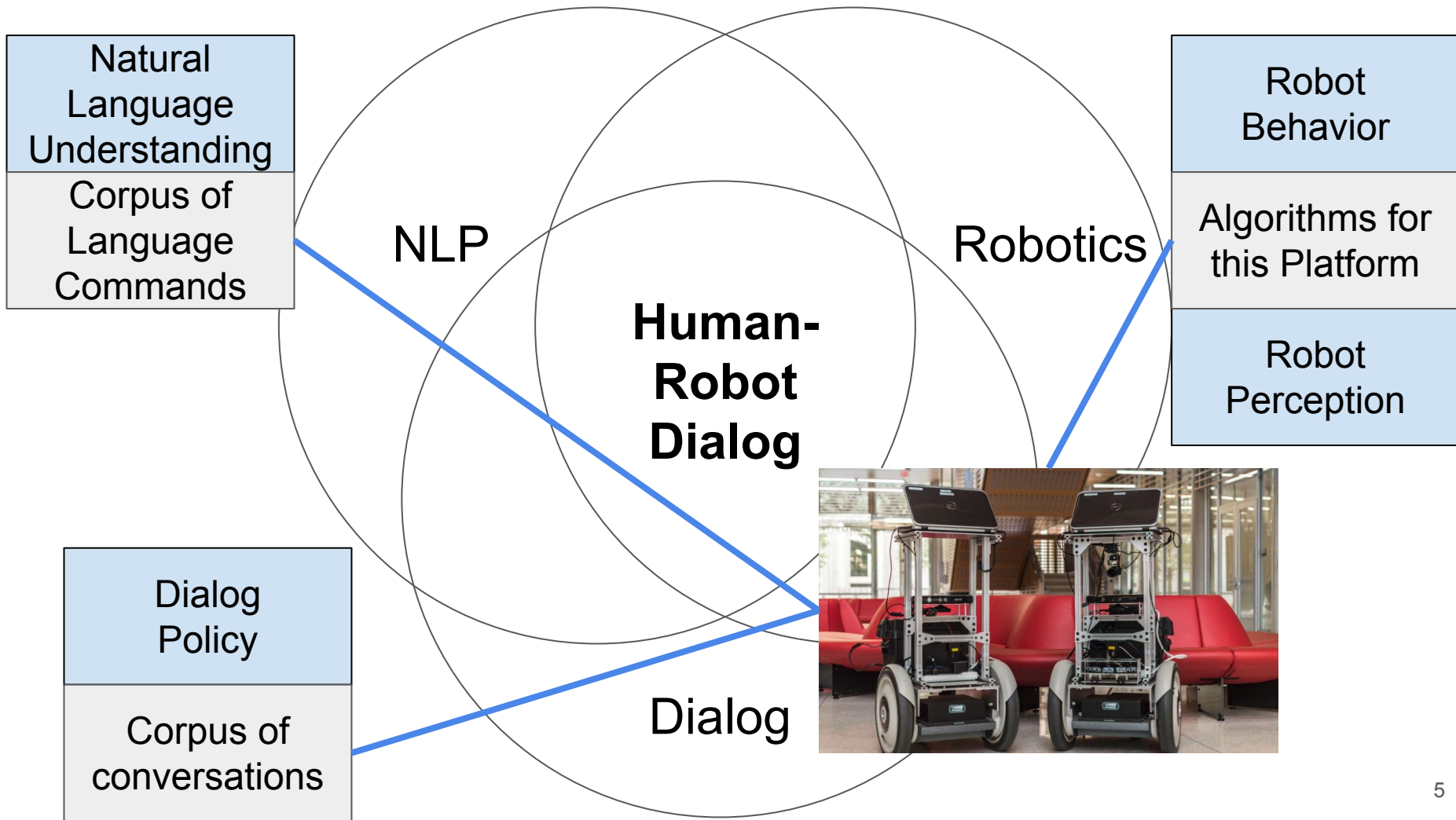
# Human-Robot Dialog

# Human-Robot Dialog



"alert me if her heart rate decreases"
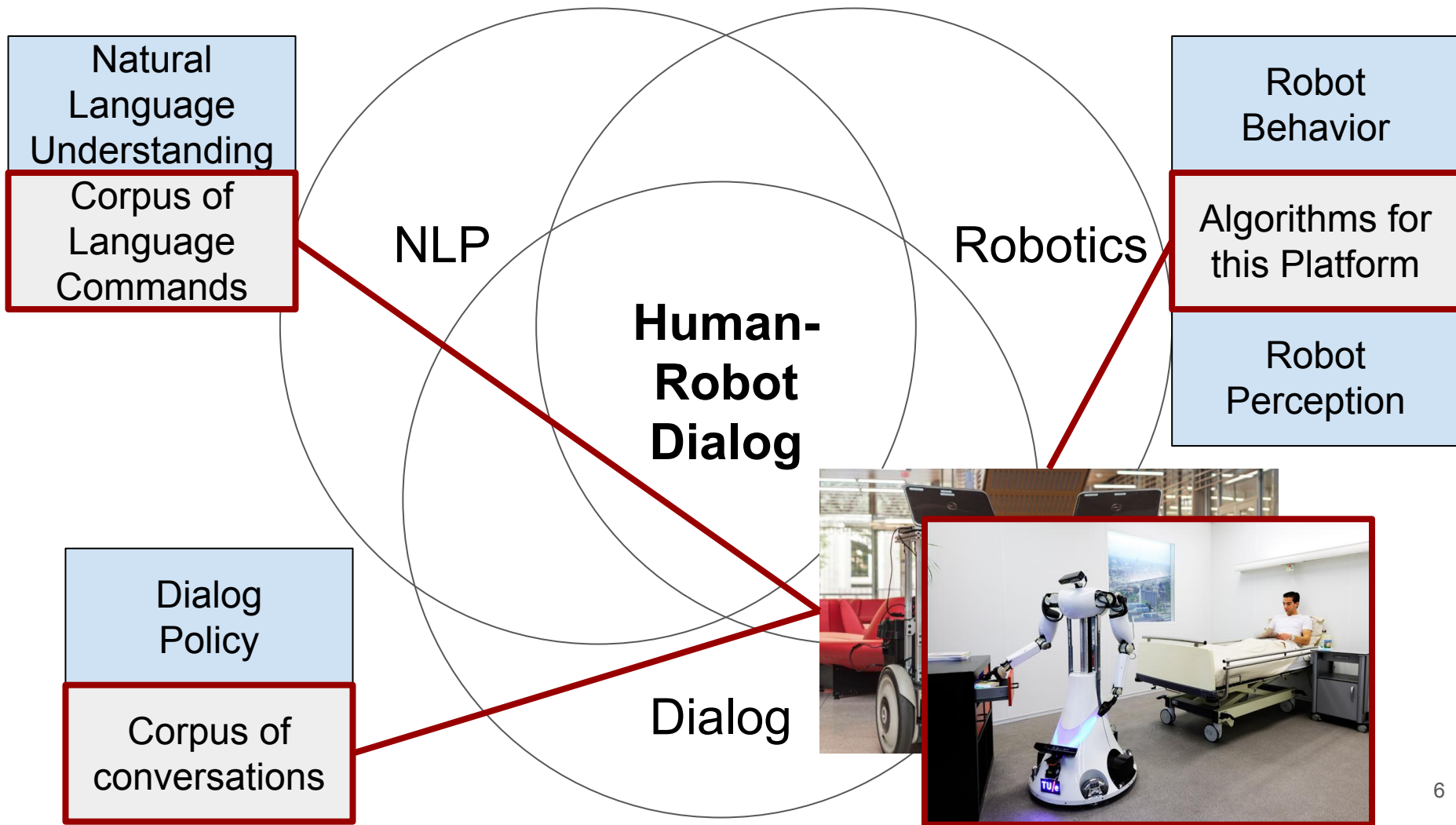"bring me his chart"
"go and get the family"
"scalpel"



"text me when the speaker arrives"
"grab the empty, green bottle"
"lead him to alice's office"
"get out of the way"

NLP

Robotics

**Human-Robot Dialog**

Dialog

Natural Language Understanding

Corpus of Language Commands

Robot Behavior

Algorithms for this Platform

Robot Perception

NLP

Robotics

**Human-Robot Dialog**

Dialog Policy

Corpus of conversations

Dialog

Natural Language Understanding

Corpus of Language Commands

NLP

Robot Behavior

Algorithms for this Platform

Robotics

Robot Perception

**Human-Robot Dialog**

Dialog Policy
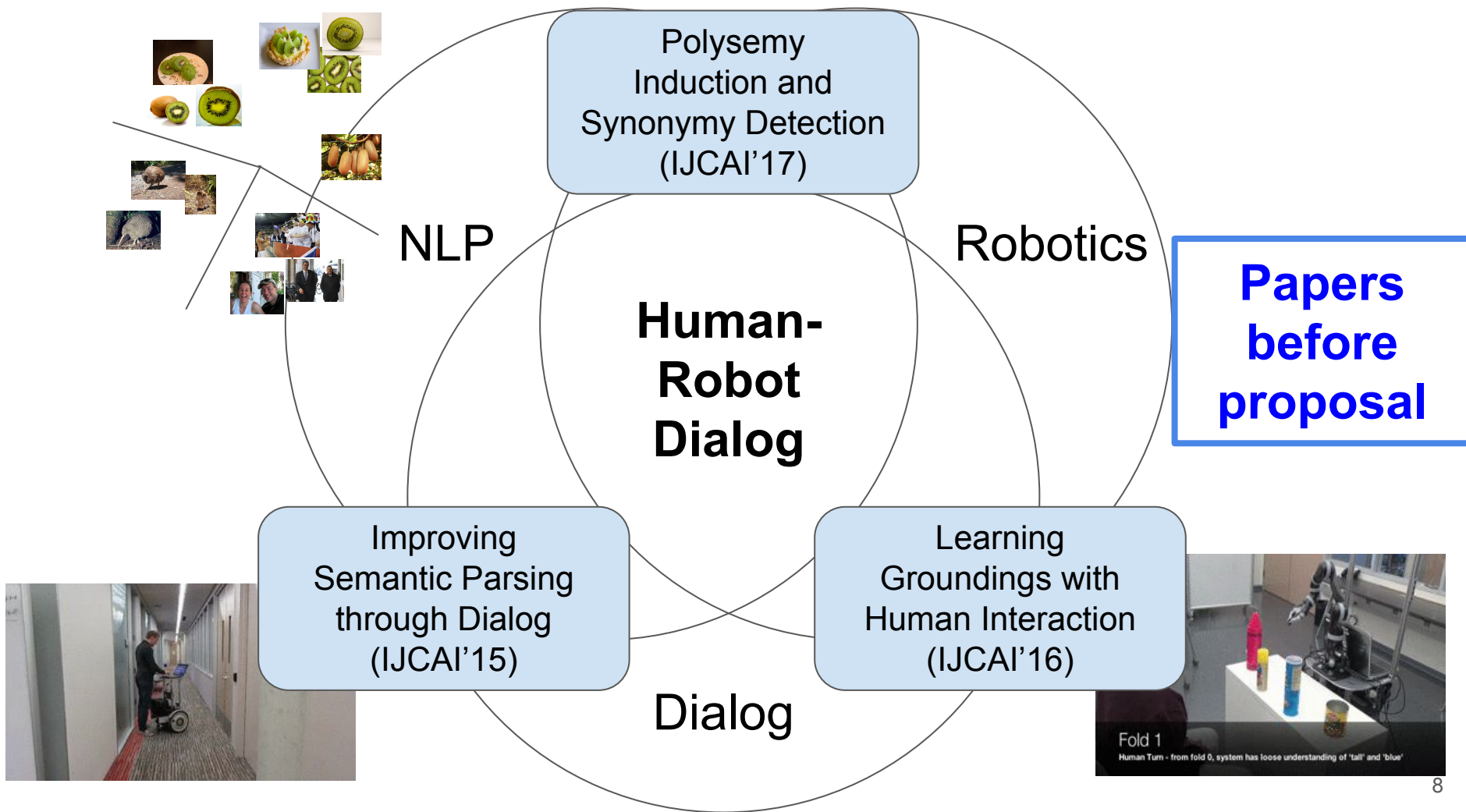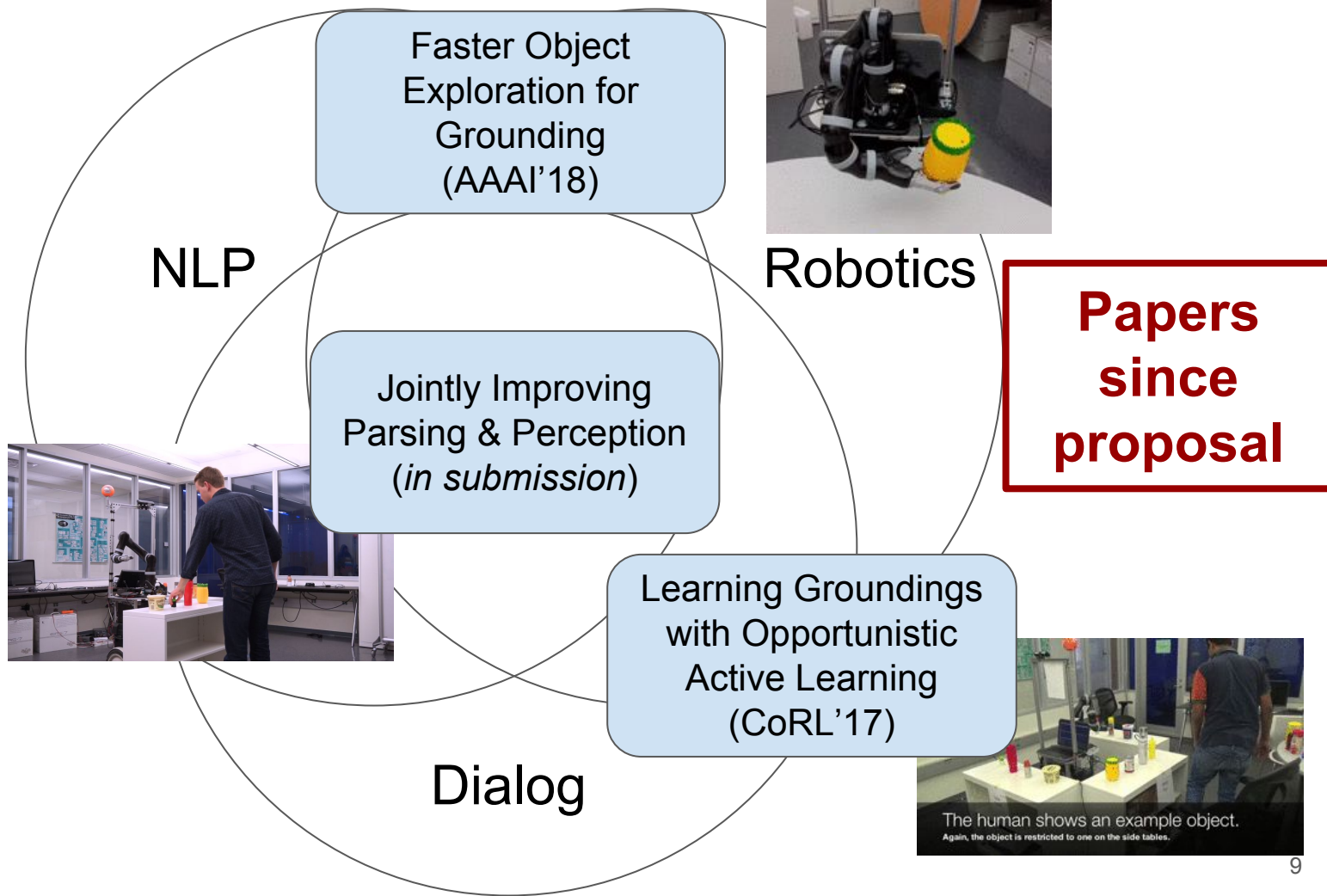
Corpus of conversations

Dialog

# Robot Dialog has Multiple Low-Resource Problems

- **My work**:

  - Develop algorithms for human-robot understanding that **overcome sparse training data.**

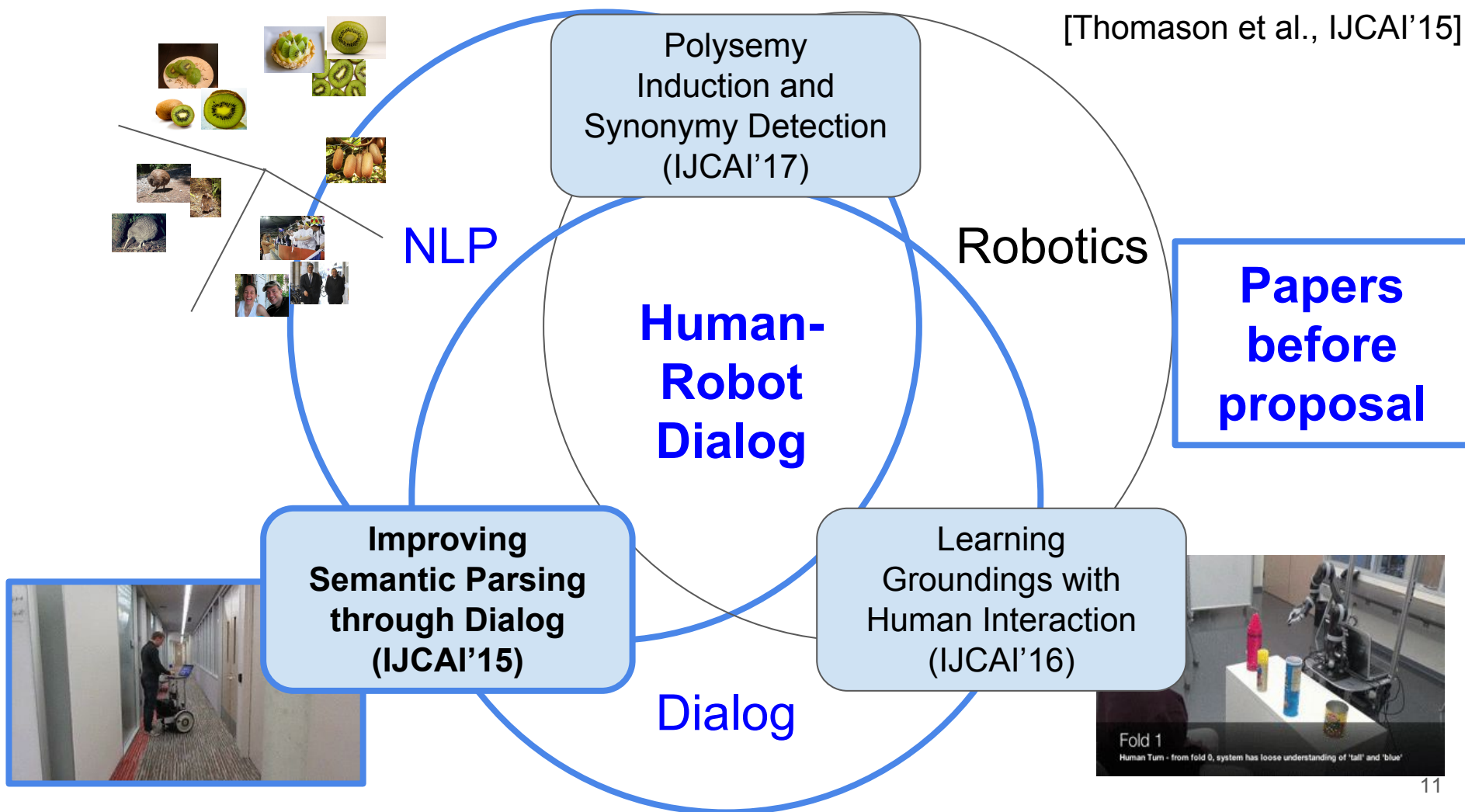  - Use dialog to **correctly perform** user requests and **better understand** future requests.

Polysemy Induction and Synonymy Detection (IJCAI'17)

NLP

Robotics

Human-Robot Dialog

Papers before proposal

Improving Semantic Parsing through Dialog (IJCAI'15)

Learning Groundings with Human Interaction (IJCAI'16)

Dialog

Fold 1
Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

NLP

Robotics

Dialog

Faster Object Exploration for Grounding (AAAI'18)

Jointly Improving Parsing & Perception (*in submission*)

Learning Groundings with Opportunistic Active Learning (CoRL'17)

**Papers since proposal**

The human shows an example object.
Again, the object is restricted to one on the side tables.

NLP

Robotics

**Human-Robot Dialog**

**Next Directions**

Dialog

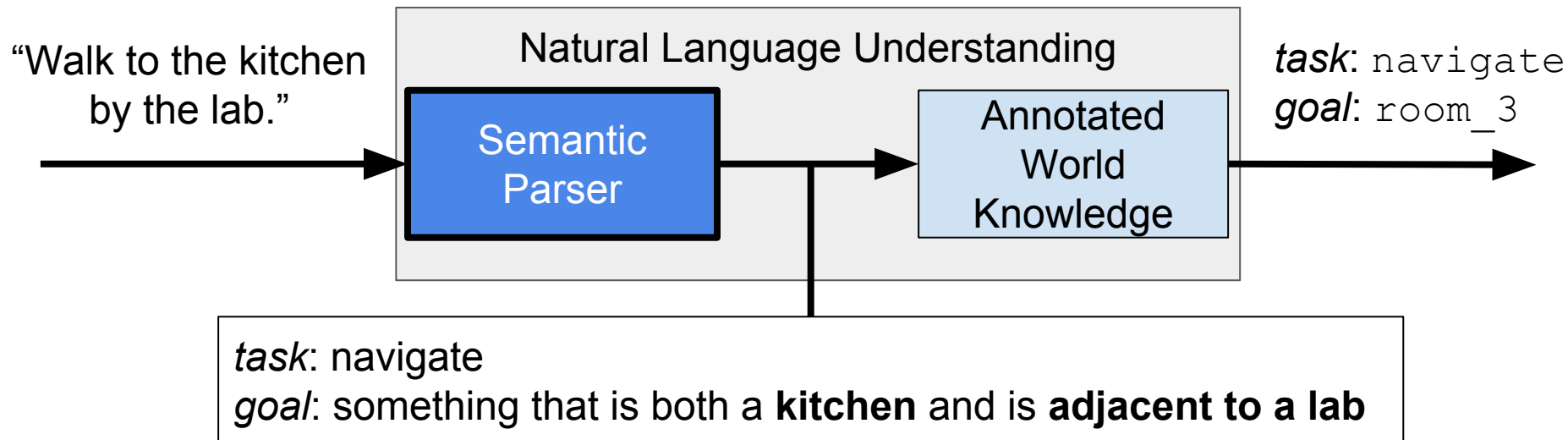[Thomason et al., IJCAI'15]

Polysemy Induction and Synonymy Detection (IJCAI'17)

NLP

Robotics

Human-Robot Dialog

Papers before proposal

Improving Semantic Parsing through Dialog (IJCAI'15)

Learning Groundings with Human Interaction (IJCAI'16)

Dialog

Fold 1
Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

11

# Dialog for Robots

User

"Walk to the kitchen by the lab." → Natural Language Understanding

*task*: `navigate`
*goal*: `room_3`

"You want me to go to room 3?"

Dialog Policy

I think I should navigate to room 3 ← Dialog Agent

# Dialog for Robots

User

"Yes."

Natural Language Understanding

*task*: `navigate`
*goal*: `room_3`

Robot Behavior

Dialog Policy

I should navigate to room 3

Dialog Agent

# Natural Language Understanding



"Walk to the kitchen by the lab."

**Natural Language Understanding**

Semantic Parser

Annotated World Knowledge

*task*: navigate
*goal*: room_3

*task*: navigate
*goal*: something that is both a **kitchen** and is **adjacent to a lab**

# Semantic Parser

- The parser can be initialized with low annotator effort.

  - In our experiments, we annotate **five** sentences.

  - Satisfies the low-resource constraints of human-robot dialog.

- But, **more training data** will improve performance.

# Inducing New Training Examples from Dialog

**TASK TO COMPLETE**

Dave Daniel wants the item in slot 5.

| | |
|---|---|
| ROBOT | How can I help? |
| YOU | please bring the item in slot 5 to dave daniel |
| ROBOT | What should I bring to dave daniel? |
| YOU | a calander |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | a day planner |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | calendar |
| ROBOT | You want me to bring calendar to dave daniel? |
| YOU | yes |
| ROBOT | I thought so |
| ROBOT | Happy to help |
| YOU | |

Final Task

Items available to robot:



16

# Inducing New Training Examples from Dialog

**TASK TO COMPLETE**
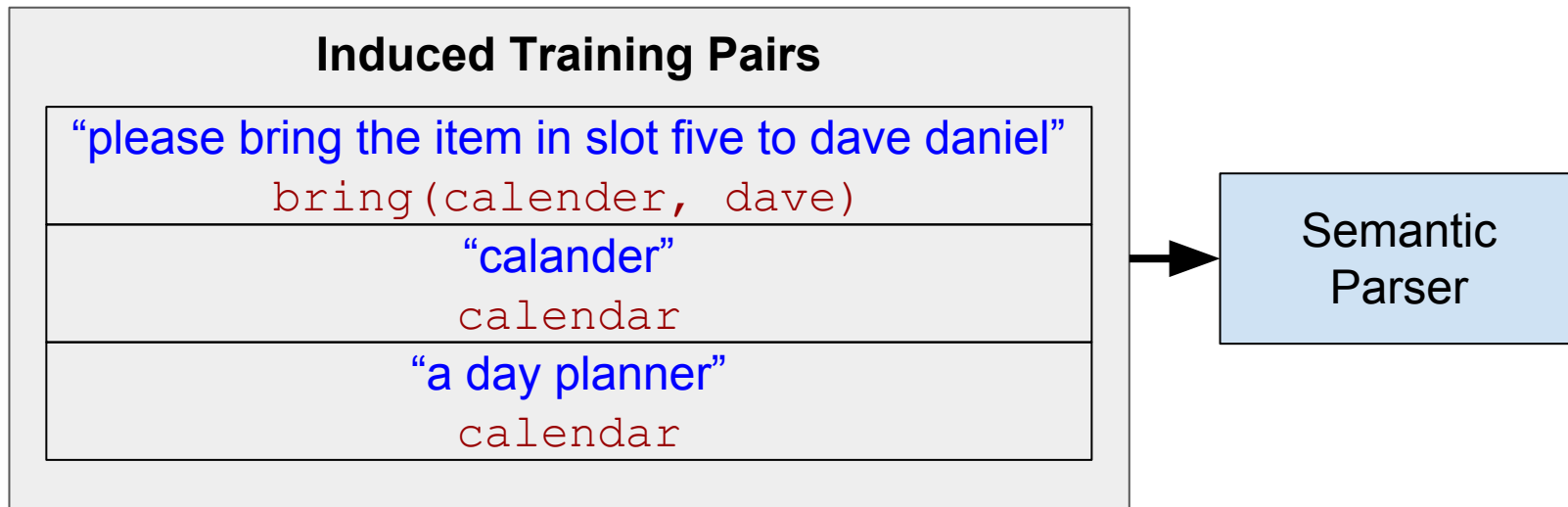
Dave Daniel wants the item in slot 5.

Items available to robot:

| ROBOT | How can I help? |
| YOU | please bring the item in slot 5 to dave daniel |
| ROBOT | What should I bring to dave daniel? |
| YOU | a calander |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | a day planner |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | calendar |
| ROBOT | You want me to bring calendar to dave daniel? |
| YOU | yes |
| ROBOT | I thought so |
| ROBOT | Happy to help |
| YOU | |

Final Task



17

# Inducing New Training Examples from Dialog

**TASK TO COMPLETE**

Dave Daniel wants the item in slot 5.

Items available to robot:

| ROBOT | How can I help? |
| YOU | please bring the item in slot 5 to dave daniel |
| ROBOT | What should I bring to dave daniel? |
| YOU | a calander |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | a day planner |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | calendar |
| ROBOT | You want me to bring calendar to dave daniel? |
| YOU | yes |
| ROBOT | I thought so |
| ROBOT | Happy to help |
| YOU | |

[Final Task]

18

# Inducing New Training Examples from Dialog

**Induced Training Pairs**

"please bring the item in slot five to dave daniel"
bring(calender, dave)

"calander"
calendar

"a day planner"
calendar

Semantic Parser

# Demonstration

# Demonstration

# Demonstration

# Demonstration

# Dialogs that Clarify Meaning and Provide Supervision

| Agent Belief (*task, goal, item, person*) | Request | Question |
|---|---|---|
| `(?, ?, ?, ?)` | *all* | "How can I help?" / "Can you reword your original request?" |
| `(navigate, ?, _, _)` | *goal* | "Where should I walk?" |
| `(deliver, _, ?, p)` | *item* | "What should I bring to `p`?" |
| `(navigate, r, _, _)` | *confirm* | "You want me to walk to `r`?" |
| . . . | . . . | . . . |

# Dialogs that Clarify Meaning and Provide Supervision

| Agent Belief (*task, goal, item, person*) | Request | Question |
|---|---|---|
| `(?, ?, ?, ?)` | *all* | "How can I help?" / "Can you reword your original request?" |
| `(navigate, ?, _, _)` | *goal* | "Where should I walk?" |
| `(deliver, _, ?, p)` | *item* | "What should I bring to `p`?" |
| `(navigate, r, _, _)` | *confirm* | "You want me to walk to `r`?" |
| . . . | . . . | . . . |

# Dialogs that Clarify Meaning and Provide Supervision

**TASK TO COMPLETE**

Dave Daniel wants the item in slot 5.

| | |
|---|---|
| ROBOT | How can I help? |
| YOU | please bring the item in slot 5 to dave daniel |
| ROBOT | What should I bring to dave daniel? |
| YOU | a calander |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | a day planner |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | calendar |
| ROBOT | You want me to bring calendar to dave daniel? |
| YOU | yes |
| ROBOT | I thought so |
| ROBOT | Happy to help |
| YOU | |

Final Task

Expect whole command

Expect *item*

Expect *item*

Expect *item*

*task*: `deliver`
*item*: `calendar`
*person:* `dave_daniel`

26

# Technical Contributions

- Design a dialog policy that allows us to **pair human language with latent meaning representations**.

- Improve semantic parsing **given very little initial in-domain data**.



**TASK TO COMPLETE**

Dave Daniel wants the item in slot 5.

| ROBOT | How can I help? |
| YOU | please bring the item in slot 5 to dave daniel |
| ROBOT | What should I bring to dave daniel? |
| YOU | a calander |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | a day planner |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | calander |
| ROBOT | You want me to bring calendar to dave daniel? |
| YOU | yes |
| ROBOT | I thought so |
| ROBOT | Happy to help |
| YOU | |

Final Task

# Experiments via Amazon Mechanical Turk

**TASK TO COMPLETE**

Dave Daniel wants the item in slot 5.

| | |
|---|---|
| ROBOT | How can I help? |
| YOU | please bring the item in slot 5 to dave daniel |
| ROBOT | What should I bring to dave daniel? |
| YOU | a calander |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | a day planner |
| ROBOT | I'm sorry, but I couldn't pinpoint what you meant by that. |
| ROBOT | What should I bring to dave daniel? |
| YOU | calendar |
| ROBOT | You want me to bring calendar to dave daniel? |
| YOU | yes |
| ROBOT | I thought so |
| ROBOT | Happy to help |
| YOU | |

Final Task

x 50

Induced Training Pairs

Semantic Parser

x 4

28

# Navigation Dialog Turns

**Navigation task
average Turker
Turns for success**

# Navigation Dialog Turns

**Navigation task average Turker Turns for success**



Robot: How can I help?
Human: go
…
Human: go to dave daniel's office

**Induced Training Pairs**

"go"
go(room_2)

...

# Delivery Dialog Turns

**Delivery task average Turker turns for success**



- Statistically significant decrease.

- More arguments:

  harder to understand, so more to

  gain from parser training.

**Qualitative**: One user wrote "the robot even fixed my typo when I mispelled calendar!"

31

# Other Findings



- Users rate system more understanding and less frustrating.

- Results replicable on physical platform.

[Thomason et al., IJCAI'15]

Polysemy Induction and Synonymy Detection (IJCAI'17)

NLP

Robotics

**Papers before proposal**

**Human-Robot Dialog**

**Improving Semantic Parsing through Dialog (IJCAI'15)**

Learning Groundings with Human Interaction (IJCAI'16)

Dialog

Fold 1
Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

33

[Thomason et al., IJCAI'16]

NLP

Robotics

Polysemy Induction and Synonymy Detection (IJCAI'17)

Human-Robot Dialog

Improving Semantic Parsing through Dialog (IJCAI'15)

Learning Groundings with Human Interaction (IJCAI'16)

Dialog

**Papers before proposal**

Fold 1

Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

34

# We do not yet handle perception information



User

"Get the empty bottle."

Natural Language Understanding

Semantic Parser

Annotated World Knowledge

Robot Behavior

Question

Meaning

Dialog Policy

Agent Belief

Dialog Agent

# We need to perform *language grounding*



"Get the empty bottle."

Natural Language Understanding

Perception Models

Semantic Parser

Annotated World Knowledge

User

Robot Behavior

Question

Dialog Policy

Meaning

Agent Belief

Dialog Agent

# Language Grounding



empty?

↓

Perception Models

↓

yes

# Language Grounding



- *Symbol grounding problem.*

- Historically use visual space.

- We use more than vision.

# Language Grounding



**_Haptic_ sensors from arm give force information.**

**_Audio_ signals from mic give sound information.**

# Perceptual Grounding

Grasp

Lift

Lower

Look

Drop

Press

Push

*color*, *shape*, and deep *VGG* features.

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|

$p$: squishy

$c$: press haptic



Few labeled examples, but SVMs can operate on this sparse data.

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|

$$\boxed{d(p,o)} = sgn\left(\sum_{c \in C} w_{p,c} \mathbf{G}_{p,c}(o)\right)$$

Decision

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|

$$d(p,o) = sgn\left(\sum_{c \in C} w_{p,c}\mathbf{G}_{p,c}(o)\right)$$

Decision    Sensorimotor Contexts

# Building Perceptual Classifiers

| | |
|---|---|
| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |

$$d(p,o) = sgn\left(\sum_{c \in C} w_{p,c} \mathbf{G}_{p,c}(o)\right)$$

Decision    Sensorimotor    Context
Contexts    SVM result

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|

$$d(p,o) = sgn\left(\sum_{c\in C} w_{p,c}\, \mathbf{G}_{p,c}(o)\right)$$

Decision    Sensorimotor   Reliability    Context

Contexts     Weight    SVM result

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|

$$d(p, o) = sgn \left( \sum_{c \in C} \boxed{w_{p,c}} \mathbf{G}_{p,c}(o) \right)$$

**Reliability weights estimated from xval**

|  | squishy |
|---|---|
| **sensorimotor context** | $w_{p,c}$ |
| press-haptics | 0.5 |
| grasp-haptics | 0.3 |
| ... | ... |
| look-VGG | 0.01 |

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|



press
haptic

Reliability weights
estimated from xval

| squishy | |
|---|---|
| **sensorimotor context** | $w_{p,c}$ |
| press-haptics | 0.5 |
| grasp-haptics | 0.3 |
| ... | ... |
| look-VGG | 0.01 |

47

# Building Perceptual Classifiers

| $\mathbf{G}_{p,c}(o)$ | SVM trained for predicate $p$ and sensorimotor context $c$ result on object $o$ |
|---|---|



look VGG

Reliability weights estimated from xval

| squishy | |
|---|---|
| **sensorimotor context** | $w_{p,c}$ |
| press-haptics | 0.5 |
| grasp-haptics | 0.3 |
| ... | ... |
| look-VGG | 0.01 |

48

# Technical Contributions

- Ensemble SVMs over **multi-modal object features** to perform **language grounding**.

- Get language labels from natural **language game** with human users

## Human Turn

Initially, the robot has no training data and randomly guesses objects.

# Experiments Playing *I Spy*



Human Turn
The description offered by the subject provides positive labels for chosen object.

Robot Turn
A follow-up dialog gives additional positive/negative labels for predicates.

**multi-modal**       vs       **vision only**

# Experiments Playing *I Spy*



Four folds of objects for
four rounds of training.



**Bold**: Lower than fold 0 average. *: Lower than vision only baseline

# Problematic *I Spy* Object



Bold: Lower than fold 0 average. *: Lower than vision only baseline

**Future**: Be mindful of object *novelty* both for the learning algorithm and for human users.

Polysemy Induction and Synonymy Detection (IJCAI'17)

NLP

Robotics

**Human-Robot Dialog**

**Papers before proposal**

Improving Semantic Parsing through Dialog (IJCAI'15)

**Learning Groundings with Human Interaction (IJCAI'16)**

Dialog

Fold 1

Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

NLP

Robotics

**Polysemy Induction and Synonymy Detection**

**Human-Robot Dialog**

**Papers before proposal**

Improving Semantic Parsing through Dialog (IJCAI'15)

Learning Groundings with Human Interaction (IJCAI'16)

Dialog

Fold 1
Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

# Unsupervised Word Synset Induction

"chinese grapefruit"

"kiwi"

"kiwi vine"

# Unsupervised Word Synset Induction

"kiwi", "chinese grapefruit",
"kiwi vine"

"kiwi"

"kiwi"

Polysemy
Induction and
Synonymy
Detection

NLP

Robotics

Human-
Robot
Dialog

**Papers
before
proposal**

Improving
Semantic Parsing
through Dialog
(IJCAI'15)

Learning
Groundings with
Human Interaction
(IJCAI'16)

Dialog

Fold 1
Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

Faster Object Exploration for Grounding (AAAI'18)

NLP

Robotics

**Papers since proposal**

Jointly Improving Parsing & Perception (*in submission*)

Learning Groundings with Opportunistic Active Learning (CoRL'17)

Dialog

The human shows an example object.
Again, the object is restricted to one on the side tables.

[Thomason et al., AAAI'18]

**Faster Object Exploration for Grounding (AAAI'18)**

NLP

Robotics

**Papers since proposal**

Jointly Improving Parsing & Perception (*in submission*)

Learning Groundings with Opportunistic Active Learning (CoRL'17)

Dialog

The human shows an example object.
Again, the object is restricted to one on the side tables.

# Exploratory Behaviors



grasp (22s)

lift (11.1s)

lower (10.6s)

+hold (5.7s)

drop (9.8s)

push (22s)

press (22s)

+look (0.8s)

104s to explore an object once.

520s to explore an object five times.

4.5 **hours** to fully explore 32 objects.

61

# Guiding Exploratory Behaviors

rigid:

squishy?

press
haptic

press?

look
VGG

look?

62

# Guiding Exploratory Behaviors

rigid:   squishy

press
haptic

**+**

press
haptic

look
VGG

**✗**

look
VGG

# Guiding Exploratory Behaviors

$$\text{similarity}(\text{rigid}, \text{squishy}) = cos(\theta)$$

# Shared Structure: Embeddings and Features



2D-projection of
word embeddings

2D-projection of
behavior context features

# Guiding Exploratory Behaviors using Embeddings

$$d(p, o) = sgn \left( \sum_{c \in C} w_{p,c} \mathbf{G}_{p,c}(o) \right)$$

$$w_{q,c} \approx \frac{1}{|P_q|} \sum_{p \in P_q} poscos(p, q) w_{p,c}$$

Surrogate reliability weights for new classifiers for $q$

Nearest word-embedding predicates to $q$

Reliability weights for trained neighbor classifiers $p$

# Technical Contributions

- Reduce exploration time when **learning a target new word**.

- Use word embeddings and human annotations to **guide behaviors.**

# Results



Color predicates   Weight predicates   Contents predicates

(dotted lines show standard error)

# Other Findings



- Human annotations help; "how would you tell if an object is *tall*?"
- Human annotations + word embeddings work better than either alone.

[Thomason et al., AAAI'18]

**Faster Object Exploration for Grounding (AAAI'18)**

NLP

Robotics

**Papers since proposal**

Jointly Improving Parsing & Perception (*in submission*)

Learning Groundings with Opportunistic Active Learning (CoRL'17)

Dialog

The human shows an example object.
Again, the object is restricted to one on the side tables.

[Thomason et al., CoRL'17]

Faster Object Exploration for Grounding (AAAI'18)

NLP

Robotics

**Papers since proposal**

Jointly Improving Parsing & Perception (*in submisison*)

**Learning Groundings with Opportunistic Active Learning**

Dialog

The human shows an example object.
Again, the object is restricted to one on the side tables.

# Active Learning for Perceptual Questions

$$o_{\min}(p) = \operatorname{argmin}_{o \in O_{tr}} (\kappa(p, o))$$

The object for which the predicate classifier is least sure of the predicted label.

d(bottle, ) = -0.6

d(bottle, ) = 0.8

d(bottle, ) = 0.4

d(bottle, ) = -0.2

# Active Learning for Perceptual Questions

| empty | |
|---|---|
| **sensorimotor context** | $w_{p,c}$ |
| lift-haptics | ? |
| lift-audio | ? |
| ... | ... |
| look-vgg | ? |

| bottle | |
|---|---|
| **sensorimotor context** | $w_{p,c}$ |
| look-shape | 0.6 |
| look-vgg | 0.5 |
| ... | ... |
| lower-haptics | 0.02 |

# Active Learning for Perceptual Questions

$$prob(p) = \frac{1 - \kappa(p, o_{\min}(p))}{\sum_{q \in P \setminus \{p\}} 1 - \kappa(q, o_{\min}(q))}$$

Ask for a label with probability proportional to *un*confidence in least confident training object.

$$p \in \{q : q \in P \wedge \kappa(q, o_{\min}(q)) = 0\}$$

Ask for a positive label for any predicate we have insufficient data for.

# Active Learning for Perceptual Questions

"Could you use the word bottle when describing this object?"



Ask for a label with probability proportional to *un*confidence in least confident training object.

"Can you show me something empty?"

Ask for a positive label for any predicate we have insufficient data for.

The human shows an example object.
Again, the object is restricted to one on the side tables.

# Technical Contributions

- Introduce an **opportunistic active learning** strategy for getting high-value labels.

- Show that *off-topic* questions **improve performance.**



"A full, yellow bottle."

"Would you describe this object as full?"

"Show me something red."

# Experiments with Object Identification



The human shows an example object.
Again, the object is restricted to one on the side tables.



Robot: Is this the object you had in mind when you said...

"Would you describe this object as full?"

**Baseline Agent**

vs

"Show me something red."

*Inquisitive* Agent

# Results



"Would you describe this object as full?"

## Baseline Agent

Rated less annoying.



"Show me something red."

## *Inquisitive* Agent

Correct object more often.

Rated better for real-world use.

[Thomason et al., CoRL'17]

NLP

Robotics

Dialog

Faster Object Exploration for Grounding (AAAI'18)

Jointly Improving Parsing & Perception (*in submission*)

**Learning Groundings with Opportunistic Active Learning**

**Papers since proposal**

The human shows an example object.
Again, the object is restricted to one on the side tables.

[*in submission*]



Faster Object Exploration for Grounding (AAAI'18)

NLP

Robotics

**Papers since proposal**

**Jointly Improving Parsing & Perception (*in submission*)**

Learning Groundings with Opportunistic Active Learning (CoRL'17)

Dialog

The human shows an example object.
Again, the object is restricted to one on the side tables.

81

# Human-Robot Dialog



Natural Language Understanding

Perception Models

Semantic Parser

Annotated World Knowledge

User

Utterance

Meaning

Robot Behavior

Question

Agent Belief

Dialog Policy

Dialog Agent

# Jointly Improving Parsing and Perception



"Move a rattling container from lounge by the conference room to Bob's office."

# Experiments via Amazon Mechanical Turk

| | |
|---|---|
| ROBOT | What should I do? |
| YOU | move a red can from the middle lab to the mallory's office |
| ROBOT | You want me to move an item from somewhere to somewhere (**not** give it to someone)? |
| YOU | yes |
| ROBOT | Where should I move something from on its way somewhere else? |
| YOU | the middle lab |
| ROBOT | To where should I move something from here? |
| YOU | type your response here... |

Say

## Training

x 113

Induced Training Pairs

Object / Predicate Labels

Semantic Parser

Perception Models

# Experiments via Amazon Mechanical Turk

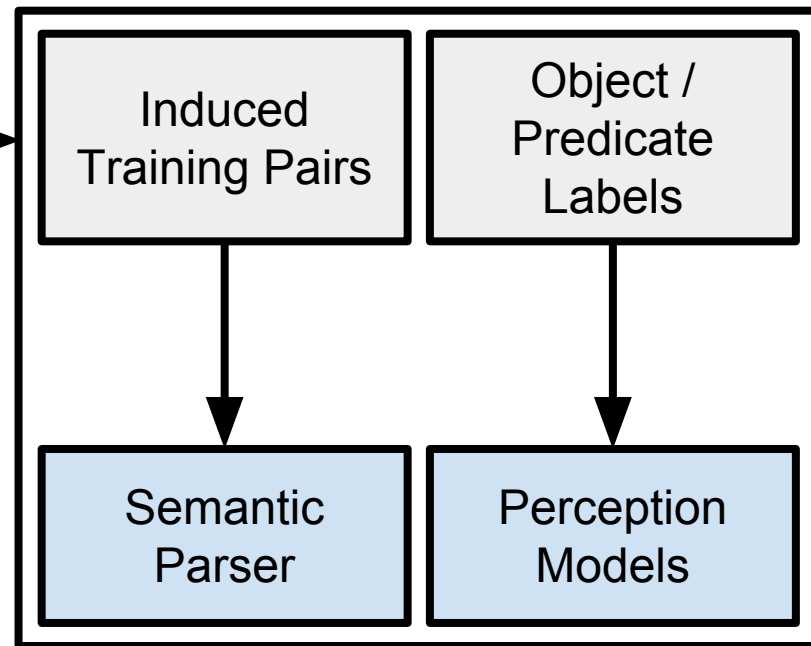| ROBOT | What should I do? |
|-------|-------------------|
| YOU | move a red can from the middle lab to the mallory's office |
| ROBOT | You want me to move an item from somewhere to somewhere (**not** give it to someone)? |
| YOU | yes |
| ROBOT | Where should I move something from on its way somewhere else? |
| YOU | the middle lab |
| ROBOT | To where should I move something from here? |
| YOU | type your response here... |

Say

x ~45

## Testing - Baseline

| Semantic Parser | Perception Models |
|-----------------|-------------------|

# Experiments via Amazon Mechanical Turk

| ROBOT | What should I do? |
|-------|-------------------|
| YOU | move a red can from the middle lab to the mallory's office |
| ROBOT | You want me to move an item from somewhere to somewhere (**not** give it to someone)? |
| YOU | yes |
| ROBOT | Where should I move something from on its way somewhere else? |
| YOU | the middle lab |
| ROBOT | To where should I move something from here? |
| YOU | type your response here... |

Say

x ~45

## Testing - Perception

Object / Predicate Labels

Semantic Parser

Perception Models

86

# Getting Object/Predicate Labels in Dialog



ROBOT | I'm still learning the meanings of some words. I'm going to ask you a few questions about these nearby objects before we continue.
ROBOT | Would you use the word 'metallic' when describing this object?
YOU | type your response here...

Say

Object / Predicate Labels

Perception Models

# Getting Object/Predicate Labels in Dialog



Object /
Predicate
Labels

Perception
Models

# Experiments via Amazon Mechanical Turk

| ROBOT | What should I do? |
|---|---|
| YOU | move a red can from the middle lab to the mallory's office |
| ROBOT | You want me to move an item from somewhere to somewhere (**not** give it to someone)? |
| YOU | yes |
| ROBOT | Where should I move something from on its way somewhere else? |
| YOU | the middle lab |
| ROBOT | To where should I move something from here? |
| YOU | type your response here... |

x ~45

## Testing - Parsing + Perception

Induced Training Pairs

Object / Predicate Labels

Semantic Parser

Perception Models

89

# Inducing New Training Examples from Dialog

| ROBOT | What should I do? |
|-------|-------------------|
| YOU | go to the middle lab |
| ROBOT | You want me to go somewhere? (**not** manipulate any objects) |
| YOU | yes |
| ROBOT | Where should I go? |
| YOU | the lab in the middle |
| ROBOT | You want me to go to here (**not** manipulate any objects)? |
| YOU | yes |
| YOU | type your response here... |

Say

The robot navigates to here.
To advance to the next task, click the button below.

Okay

Induced Training Pairs

Semantic Parser

# Inducing New Training Examples from Dialog

| | |
|---|---|
| ROBOT | What should I do? |
| YOU | go to the middle lab |
| ROBOT | You want me to go somewhere? (**not** manipulate any objects) |
| YOU | yes |
| ROBOT | Where should I go? |
| YOU | the lab in the middle |
| ROBOT | You want me to go to here (**not** manipulate any objects)? |
| YOU | yes |
| YOU | type your response here... |

Say

The robot navigates to here.
To advance to the next task, click the button below.

Okay

Expect whole command

Expect *goal*

*task*: `navigate`
*goal*: `room_3`

# Inducing New Training Examples from Dialog

**Induced Utterance/Denotation Pairs**

"go to the middle lab"
navigate(room_3)

"the lab in the middle"
room_3

# Natural Language Understanding



"the lab in the middle" → Semantic Parser → Perception Models / Annotated World Knowledge → room_3

something that is a **lab**
something that is both a **lab** and is **central**
something that is **central**
...

room_3, room_7, ...
room_3
room_3, room_1, ...
...

# Inducing New Training Examples from Dialog

# Using Embeddings for Out-of-Vocabulary Words



Induced Training Pairs

*task*: `deliver`
*item*: `coffee`
*person*: `bob`

"deliver
java to bob"

**Word Embeddings**

"deliver
java to bob"

"deliver" -> "bring"
"java" -> "coffee"

Semantic Parser

95

# Using Embeddings to Find Perception Words



| | |
|---|---|
| ROBOT | What should I deliver to this person? |
| YOU | a tall metal can |
| ROBOT | I haven't heard the word 'tall' before. Does it refer to properties of things, like a color, shape, or weight? |
| YOU | yes |
| ROBOT | Does 'tall' mean the same thing as 'tower'? |
| YOU | no |
| ROBOT | Does 'tall' mean the same thing as 'white'? |
| YOU | no |
| ROBOT | Does 'tall' mean the same thing as 'long'? |
| YOU | yes |

# Technical Contributions

● Improve **both parsing and perception** from conversations.

● Use word embeddings to **guide search for synonyms** and **novel perceptual predicates**.



97

# Experiments via Amazon Mechanical Turk

## Untrained Baseline

| Semantic Parser | Perception Models |
|---|---|

## Perception Training

Object / Predicate Labels

↓

| Semantic Parser | Perception Models |
|---|---|

## Parsing + Perception Training

| Induced Training Pairs | Object / Predicate Labels |
|---|---|

↓ ↓

| Semantic Parser | Perception Models |
|---|---|

# Metric - Semantic F1

$$T_U = \{(\text{action}, \text{deliver}), (\text{patient}, o_2), (\text{recipient}, p_1)\},$$

$$T_G = \{(\text{action}, \text{relocate}), (\text{patient}, o_2), (\text{source}, r_1), (\text{goal}, r_3)\};$$

$$\text{precision}(T_U, T_G) = \frac{|T_U \cap T_G|}{|T_U|} = \frac{1}{3},$$

$$\text{recall}(T_U, T_G) = \frac{|T_U \cap T_G|}{|T_G|} = \frac{1}{4},$$

$$f(T_U, T_G) = 2 \cdot \frac{\text{precision}(T_U, T_G) \cdot \text{recall}(T_U, T_G)}{\text{precision}(T_U, T_G) + \text{recall}(T_U, T_G)} = 0.286.$$

# Results - Navigation Task



Quantitative - Semantic F1

Qualitative - Usability Rating

# Results - Delivery Task

# Results - Relocation Task

## Quantitative - Semantic F1



## Qualitative - Usability Rating

Faster Object Exploration for Grounding (AAAI'18)

NLP

Robotics

**Papers since proposal**

**Jointly Improving Parsing & Perception (*in submission*)**

Learning Groundings with Opportunistic Active Learning (CoRL'17)

Dialog

NLP

Robotics

**Human-Robot Dialog**

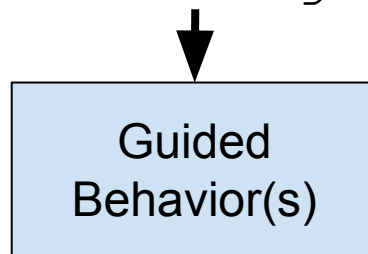**Next Directions**

Dialog

# Grounded Predicate Synset Induction



"pale"

"light"

"small"

# Grounded Predicate Synset Induction



"light"/"pale"

"light"/"small"

# Guided Exploration of New Objects



"**Move** a **rattling container** from **the kitchen** to **bob's office**."

`rattling`**?**

Guided Behavior(s)

Perception Models

**yes / no**

# Moving Forward

- The intersection of problems in human-robot dialog is **inherently low-resource**.

- Other parts of NLP, Robotics, and Dialog are not.

- We can **use big data and techniques** from these fields when solving problems in human-robot dialog.

# Moving Forward - Using Big Data Where We Can

**Very Large Corpus of Unstructured Text**

→

**Latent Language Information**
*Word Embeddings*
World Knowledge
Statistical Scripts

...

# Moving Forward - Using Big Data Where We Can

**Very Large Corpus of Training Examples**
Crowd-sourced (ImageNet)

VGG Net

good features

bottle

111

# Moving Forward - Using Big Data Where We Can



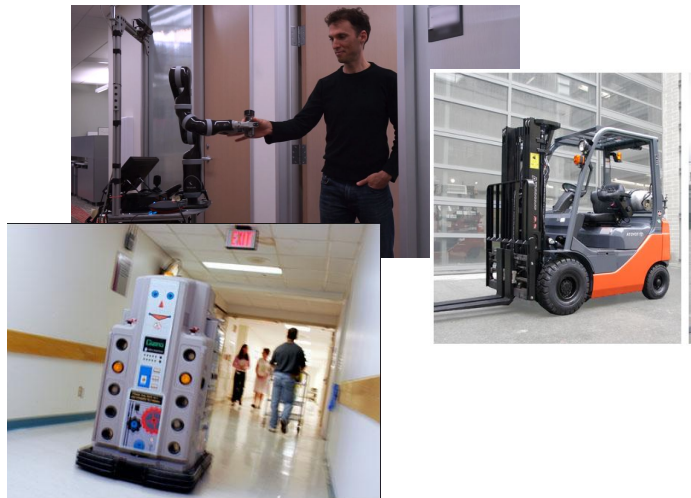**Corpus of Object Representations from Exploratory Behaviors**

**Latent Representations**
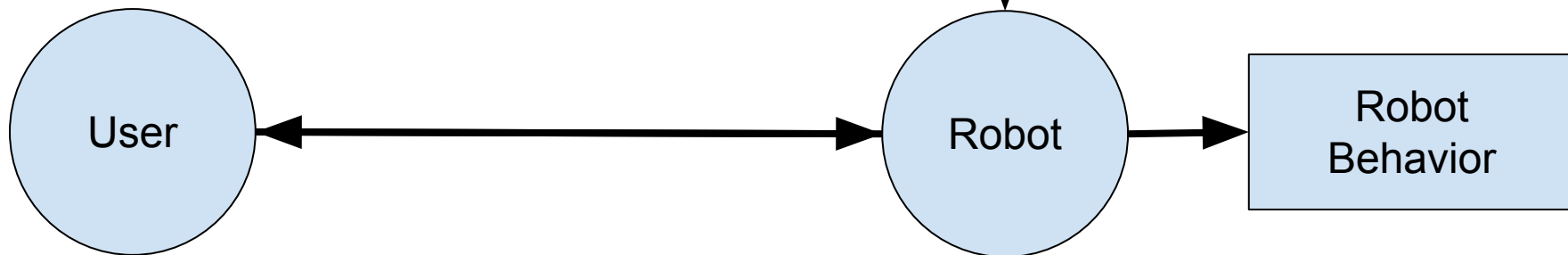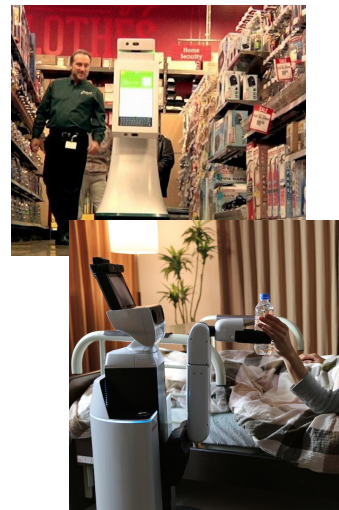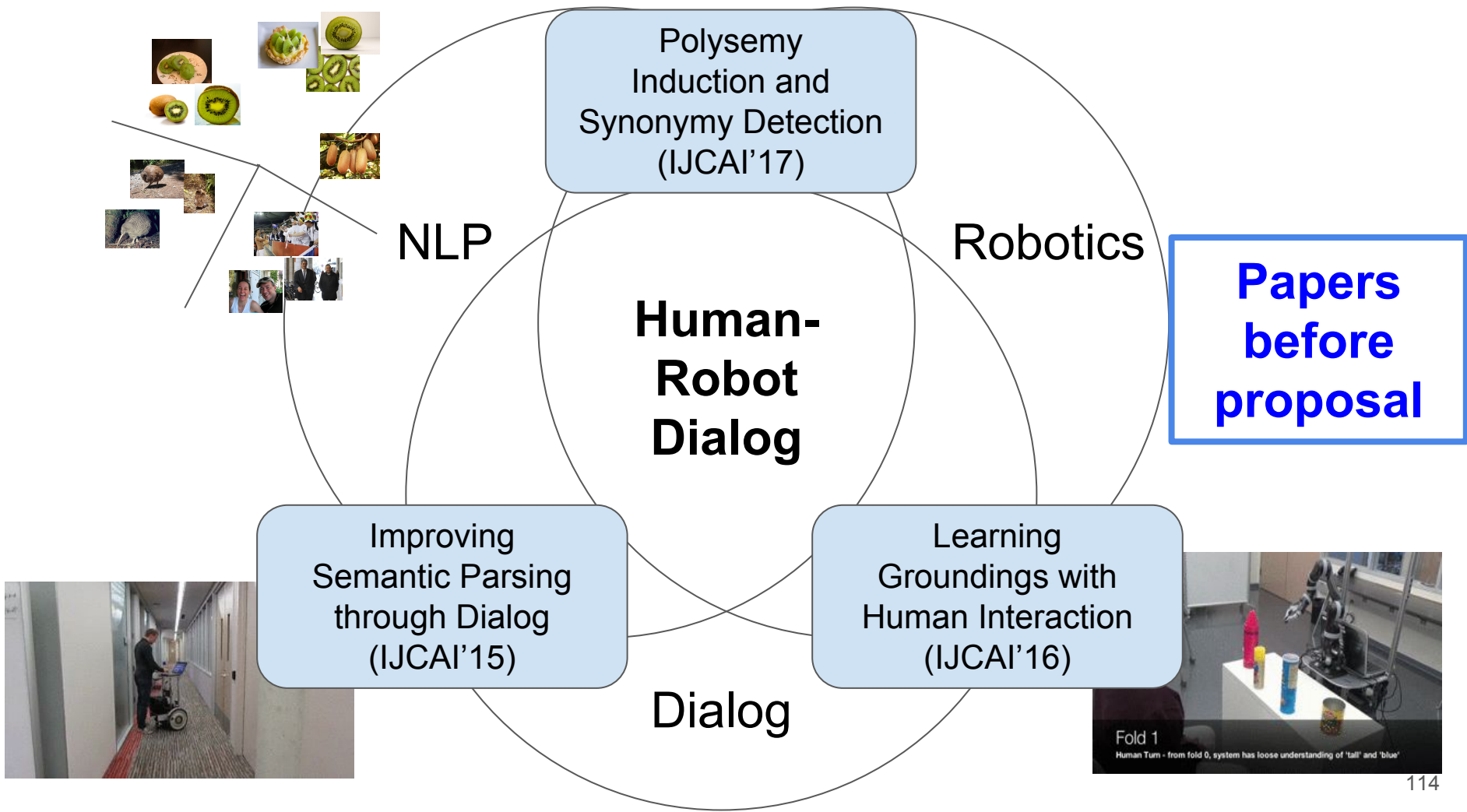Autoencoders
GANs
....

good features?
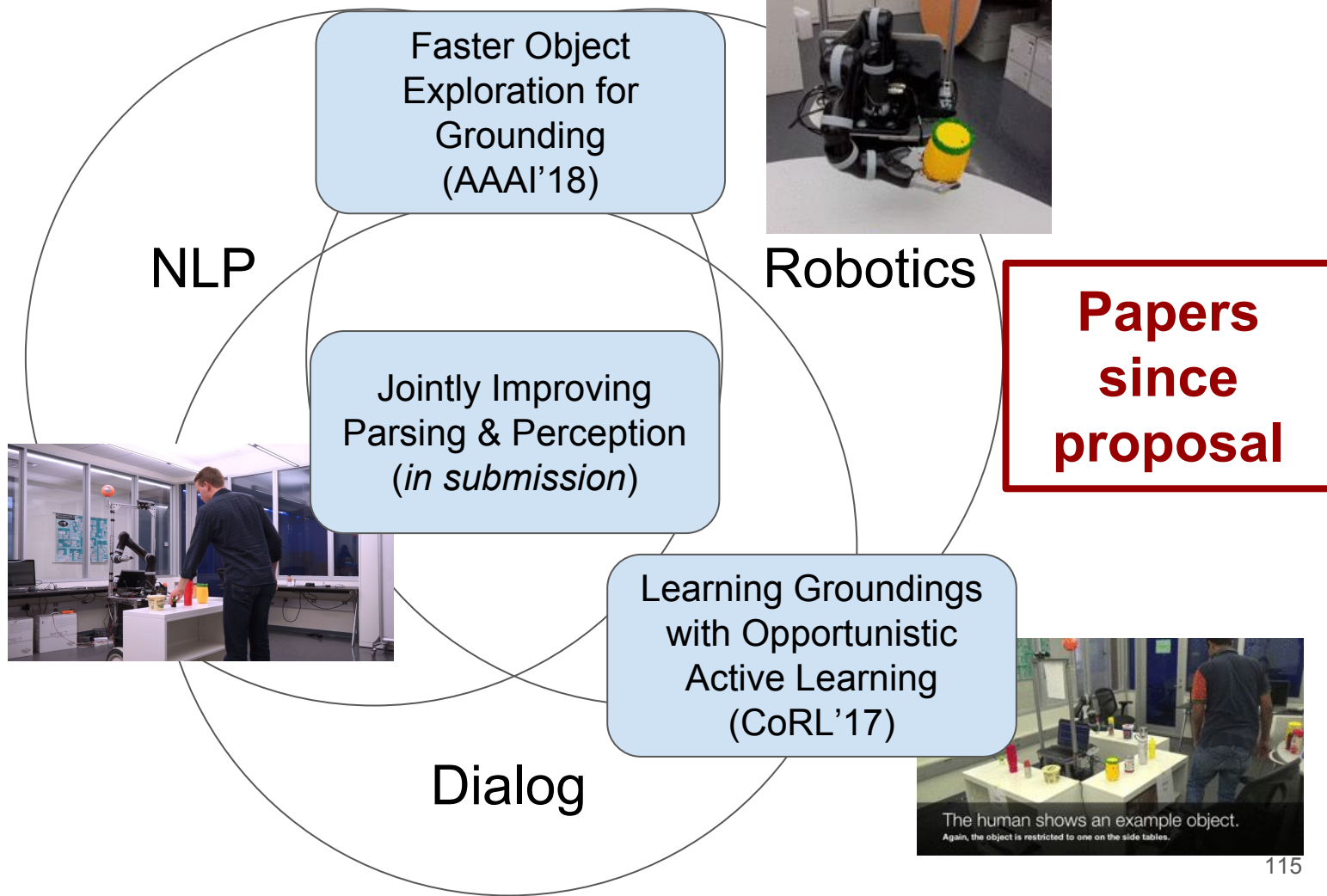
# Moving Forward - Transfer Learning



**Corpus of Human-Robot Dialogs**

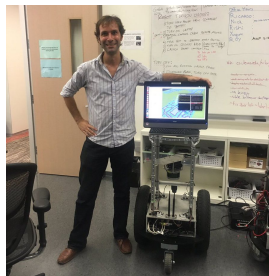Similar domain shared commands
Sharing object representations

User ↔ Robot → Robot Behavior

Polysemy Induction and Synonymy Detection (IJCAI'17)

NLP

Robotics

Human-Robot Dialog

Papers before proposal

Improving Semantic Parsing through Dialog (IJCAI'15)

Learning Groundings with Human Interaction (IJCAI'16)

Dialog

Fold 1
Human Turn - from fold 0, system has loose understanding of 'tall' and 'blue'

114

NLP

Faster Object Exploration for Grounding (AAAI'18)

Robotics

**Papers since proposal**

Jointly Improving Parsing & Perception (*in submission*)

Learning Groundings with Opportunistic Active Learning (CoRL'17)

Dialog

The human shows an example object.
Again, the object is restricted to one on the side tables.

# Acknowledgments



Ray
Mooney

Peter
Stone

Scott
Niekum

Stefanie
Tellex

# Acknowledgments



Jivko Sinapov

Shiqi Zhang

Aishwarya Padmakumar

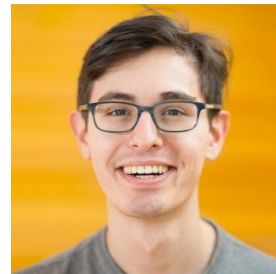Piyush Khandelwal

Rodolfo Corona

Harel Yedidsion

Justin Hart

Subhashini Venugopalan

Yuqian Jiang

Nick Walker

- *Jointly Improving Parsing and Perception for Natural Language Commands through Human-Robot Dialog*.
  **Jesse Thomason**, Aishwarya Padmakumar, Jivko Sinapov, Nick Walker, Harel Yedidsion, Justin Hart, Peter Stone, Raymond J. Mooney. (*in submission*)
- *Guiding Exploratory Behaviors for Multi-Modal Grounding of Linguistic Descriptions*.
  **Jesse Thomason**, Jivko Sinapov, Raymond J. Mooney, and Peter Stone. AAAI'18.
- *Improving Black-box Speech Recognition using Semantic Parsing*.
  Rodolfo Corona, **Jesse Thomason**, and Raymond J. Mooney. IJCNLP'17.
- *Opportunistic Active Learning for Grounding Natural Language Descriptions*.
  **Jesse Thomason**, Aishwarya Padmakumar, Jivko Sinapov, Justin Hart, Peter Stone, and Raymond J. Mooney. CoRL'17.
- *Multi-Modal Word Synset Induction*.
  **Jesse Thomason** and Raymond J. Mooney. IJCAI'17.
- *Integrated Learning of Dialog Strategies and Semantic Parsing*.
  Aishwarya Padmakumar, **Jesse Thomason**, Raymond J. Mooney. EACL'17.
- *BWIBots: A platform for bridging the gap between AI and human--robot interaction research*.
  Piyush Khandelwal, Shiqi Zhang, Jivko Sinapov, Matteo Leonetti, **Jesse Thomason**, Fangkai Yang, Ilaria Gori, Maxwell Svetlik, Priyanka Khante, Vladimir Lifschitz, J. K. Aggarwal, Raymond Mooney, and Peter Stone. IJRR'17.
- *Learning Multi-Modal Grounded Linguistic Semantics by Playing "I Spy"*.
  **Jesse Thomason**, Jivko Sinapov, Maxwell Svetlik, Peter Stone, and Raymond J. Mooney. IJCAI'16.
- *Learning to Interpret Natural Language Commands through Human-Robot Dialog*.
  **Jesse Thomason**, Shiqi Zhang, Raymond J. Mooney, and Peter Stone. IJCAI'15.

# Graded Adjectives

- Think of gradation as a form of polysemy

- Semantic parser can use surrounding context

- Re-ranking of parses, as discussed, can help disambiguate
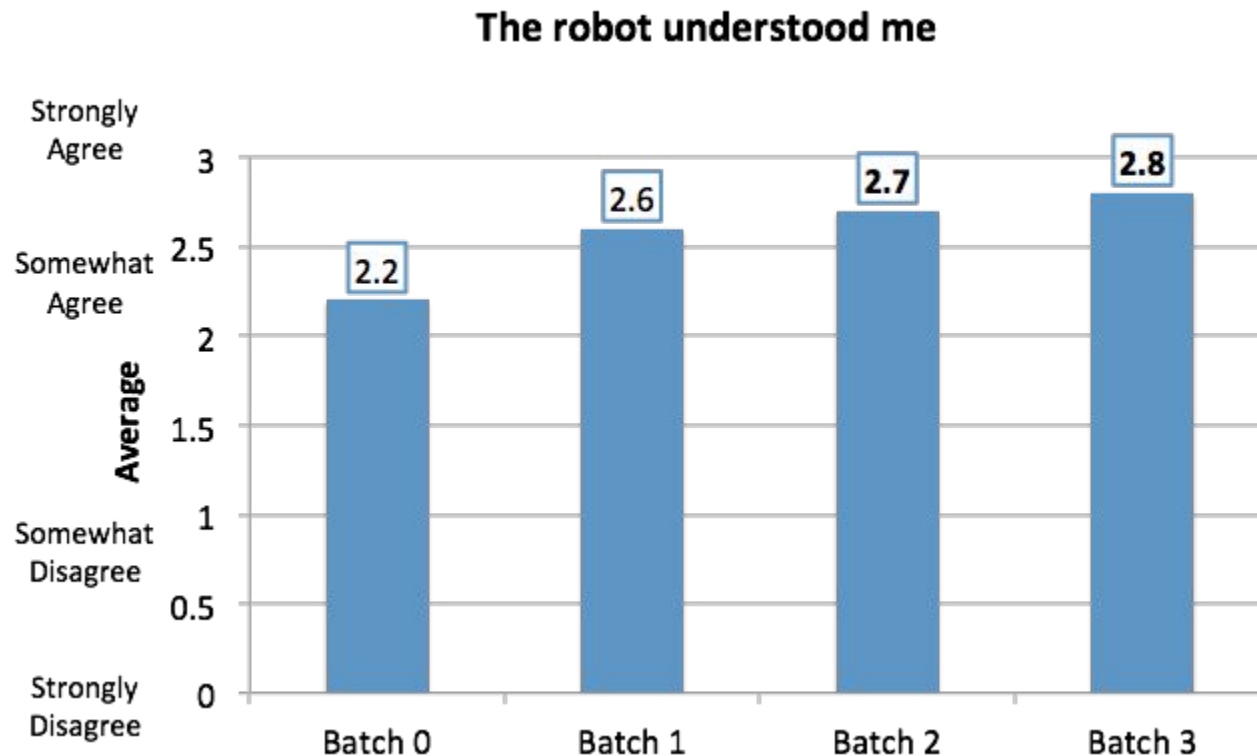
**words**



"plate"

"heavy"

"mug"

**words**

**predicates**

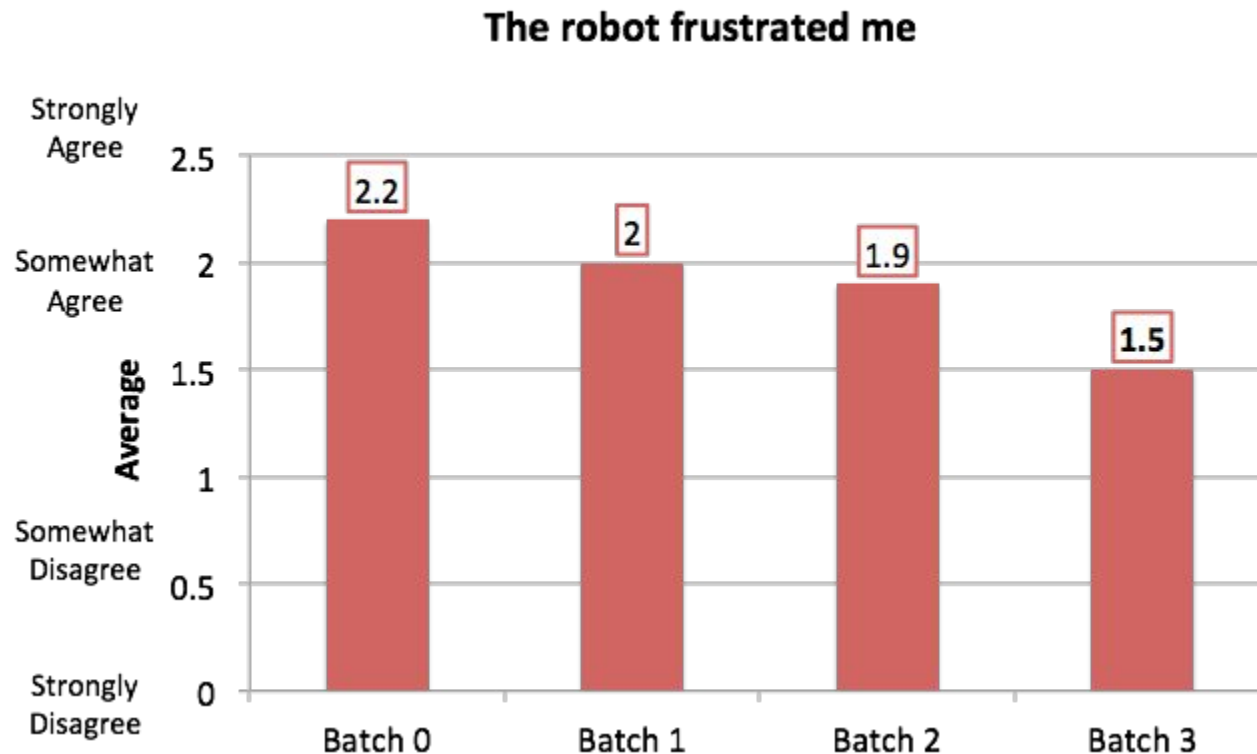"plate" — plate0

"heavy" — heavy1 / heavy0

"mug" — mug0

# Comparative Adjectives

- E.g. "taller", "heavier"; take two arguments: obj1, obj2

- Train classifier on the feature differences between obj1, obj2

- Can otherwise be handled with existing architecture

- Superlatives: majority winner object in pairwise comparative

# Mechanical Turk Qualitative Results

**The robot understood me**

# Mechanical Turk Qualitative Results

**The robot frustrated me**

# Multi-modal Representation

[Thomason et al., IJCAI'17;
Deerwester et al., 1990;
Simonyan and Zisserman, CoRR'14]

- LSA embedding text features; VGG image features

**Bat**

"... most of the oldest known, definitely identified bat fossils were already very similar to modern microbats … "

**Bat**

"... a baseball bat is divided into several regions …"

**Bat**

"... about 70% of bat species are insectivores … "

**Bat**

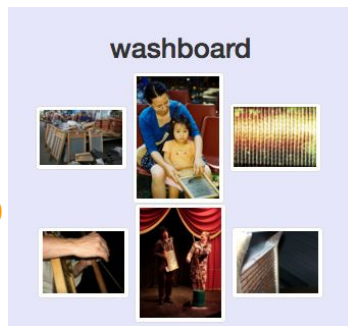"... hickory has fallen into disfavor over its greater weight, which slows down bat speed … "

# Technical Contributions

- Perform **unsupervised, multi-modal** sense induction and synonymy detection

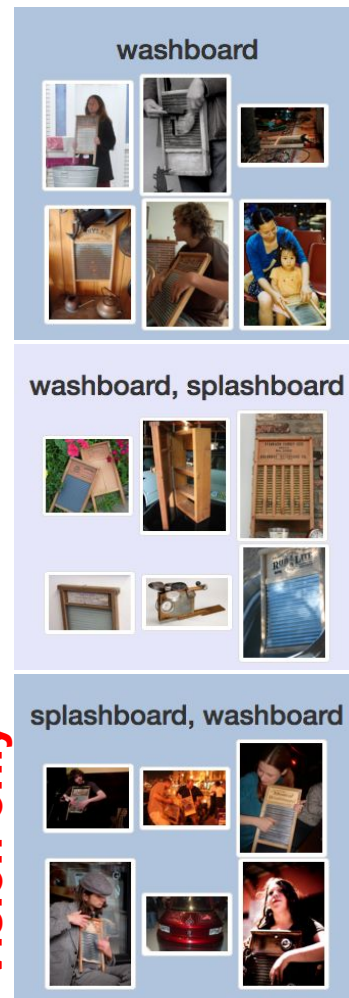- Create an ImageNet-like resource **without manual annotation**.
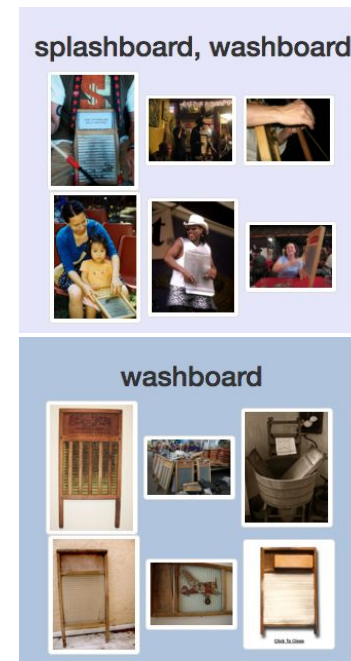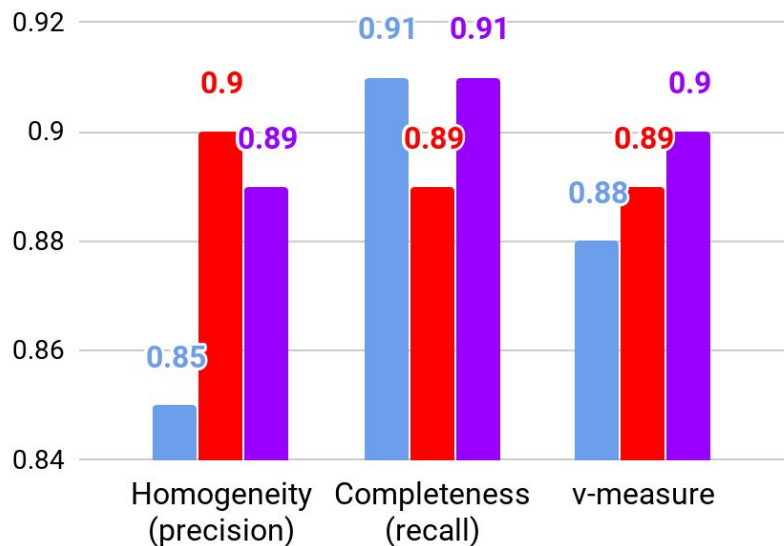
# Results

**ImageNet**

splashboard, washboard



washboard



**Text-only**

psaltery, washboard, dulcimer, cithern, headstock



king post, dugout, washboard, catapult, knothole



**Vision-only**

washboard



washboard, splashboard



splashboard, washboard



**Multi-modal**

splashboard, washboard



washboard

# Results



Synset Agreement with ImageNet

Human Evaluation

- text-only
- vision-only
- multi-modal
- ImageNet

# Results - Correct Object Selected



Correct Guess

Same Question Budget

# Results - Users Feeling Understood



The robot seemed to understand my descriptions.

Same Question Budget

# Results - Users Annoyed



The robot asked too many questions.

Baseline ■ Inquisitive    Same Question Budget

Average Likert Response vs Round

# Results - Viable for Deployment

I would use a robot like this to get objects for me in another room.

Same Question Budget



132

# Learning from Denotations

- Given utterance-denotation pair, find a semantic form that is plausible for both

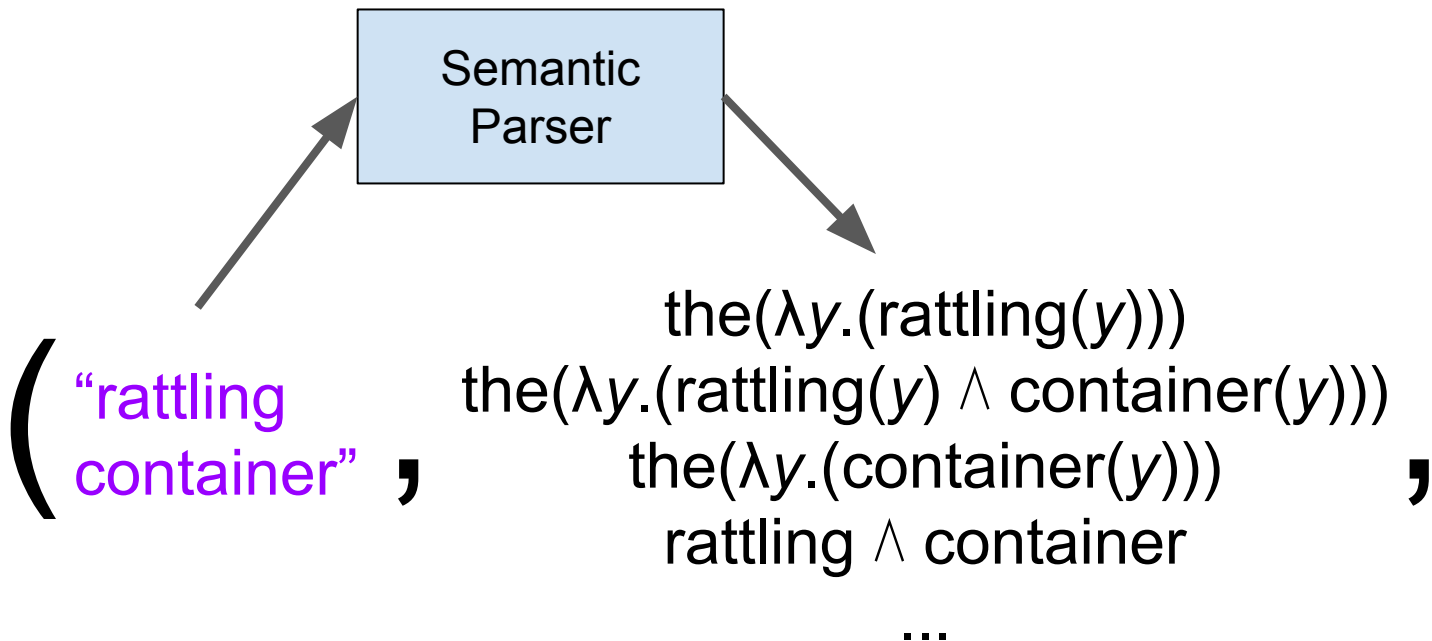( "rattling container" ,  )

# Learning from Denotations

- Use the parser to produce a beam of parses

- Use the grounder to find the denotations of those parses

# Learning from Denotations
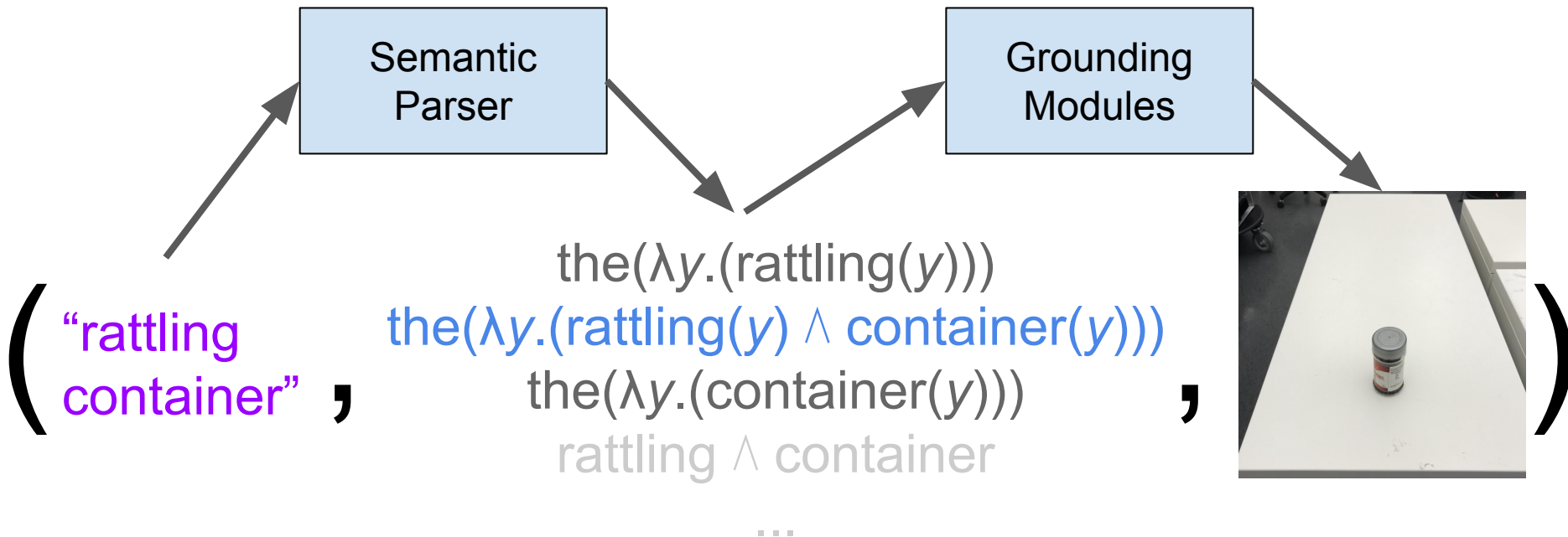


Semantic Parser

$\Big($ "rattling container" , $\quad$ the($\lambda y$.(rattling($y$)))
the($\lambda y$.(rattling($y$) $\wedge$ container($y$)))
the($\lambda y$.(container($y$)))
rattling $\wedge$ container
... , $\quad$ $\Big)$

# Learning from Denotations



Semantic
Parser

Grounding
Modules

$\Big($ "rattling container" $,$ the($\lambda y.$(rattling($y$)))
the($\lambda y.$(rattling($y$) $\wedge$ container($y$)))
the($\lambda y.$(container($y$)))
rattling $\wedge$ container

... $,$ $\Big)$

# Learning from Denotations

$\Big($ "rattling container" , the($\lambda y.$(rattling($y$) $\wedge$ container($y$))) ,  $\Big)$

# Learning from Denotations

$$\left( \text{``rattling container''} \, , \, \text{the}(\lambda y.(\text{rattling}(y) \wedge \text{container}(y))) \right)$$
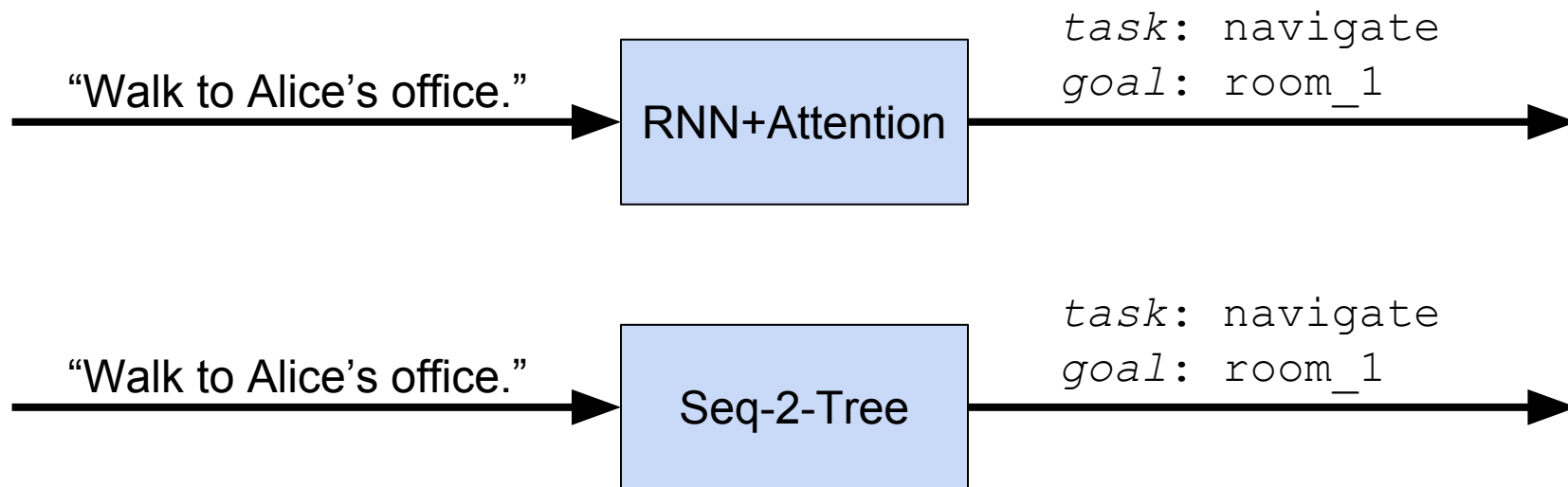
Robot
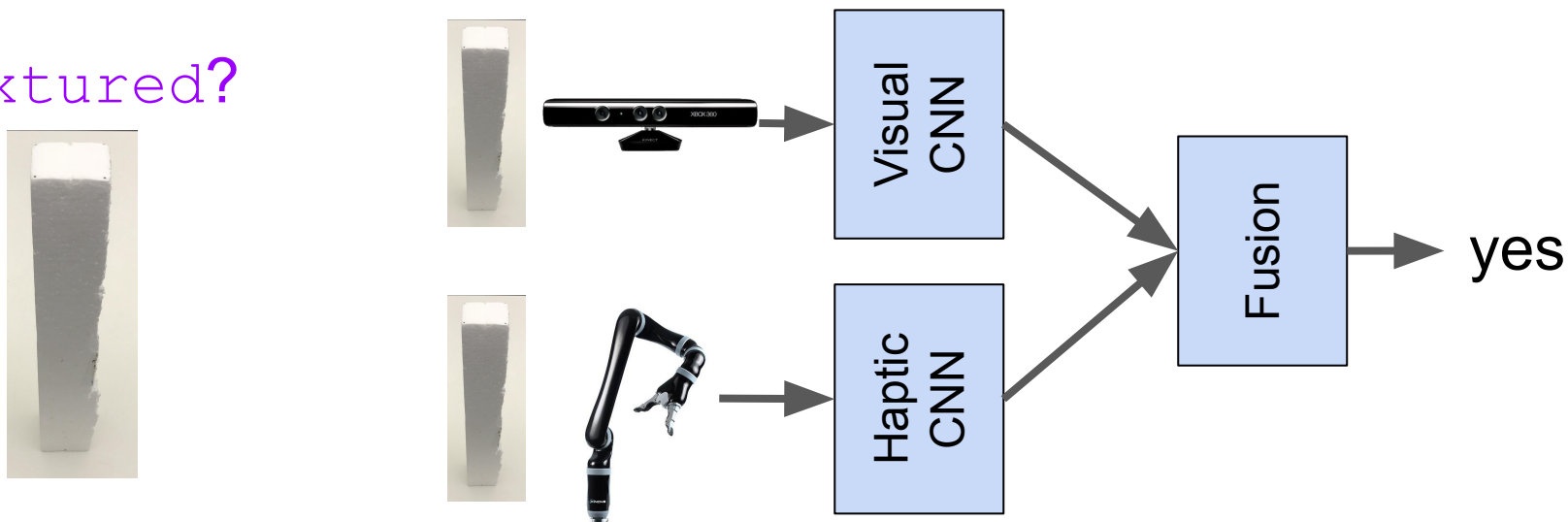"You want me to move an item from 3516 to 3510?"

# Neural Parsing Methods

- Recurrent Neural Networks (RNNs) with Attention

- Sequence-to-Tree encoder-decoder networks

"Walk to Alice's office." → **RNN+Attention** → *task*: navigate *goal*: room_1

"Walk to Alice's office." → **Seq-2-Tree** → *task*: navigate *goal*: room_1

# Neural Perception Models

- Compress high-dimensional sensorimotor context information using Convolutional Neural Networks (CNNs)



textured?

yes

# Embodied Question Answering

- End-to-end deep model for joint parsing and perception